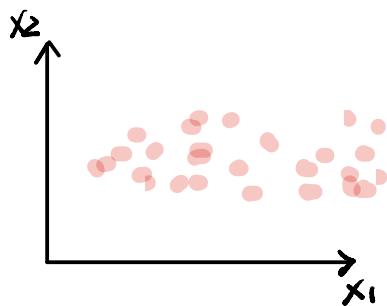
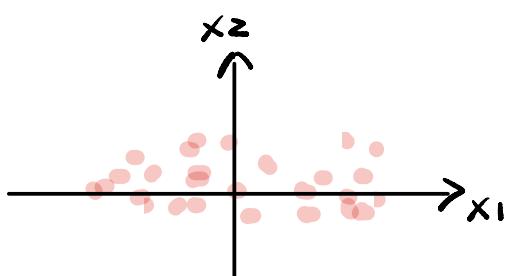


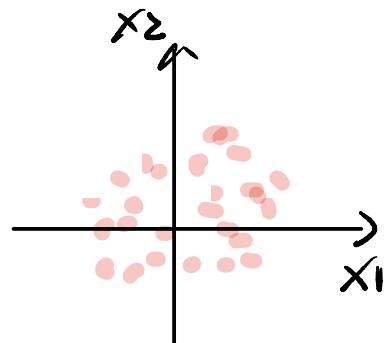
Normalizing training sets



State One



State Two



State Three.

Step one: Subtract Mean (State 1 \rightarrow State 2).

$$\mu = \frac{1}{m} \sum_{i=1}^m x^{(i)}$$

$$x := x - \mu$$

Effect: Centering data around mean.

Step two: Normalize variance (Stage 2 \rightarrow Stage 3).

$$\sigma^2 = \frac{1}{m} \sum_{i=1}^m x^{(i)} * x^{(i)}$$

$x' = \frac{x - \mu}{\sigma}$: subtract mean and divided by standard deviation.

* same σ and μ need to be applied to normalize test set.

Why normalize input?

$$w_1 = x_1 : 1 - 1000$$

$$w_2 = x_2 : 1 - 2$$

Since $J(w, b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}, y)$, and $\hat{y} = w_1 x_1 + w_2 x_2 + \dots + w_n x_n + b$, w_1 in example on the left will have bigger effect on our result than w_2 , which is an assumption that we don't want to make.

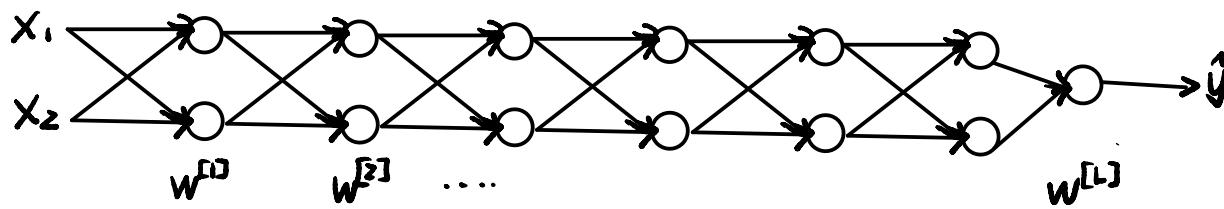
Another reason for normalizing inputs:

Neural Networks with Unnormalized inputs have more oval cost function shape comparing a more circular cost function resulting from normalized neural networks. As a result, gradient descent will do more zigzagging movements and increase training expense in unnormalized networks.

In short: Unnormalized inputs = harder to learn thru gradient descent.

Vanishing / exploding gradients.

Imagine a very deep neural network:



Assume $g(z) = z$, $b^{[l]} = 0$

$$\rightarrow \hat{y} = \underline{w^{[L]} w^{[L-1]} \dots w^{[1]} x}$$

if $w^{[l]} > I$, e.g. $w^{[L]} = \begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix}$

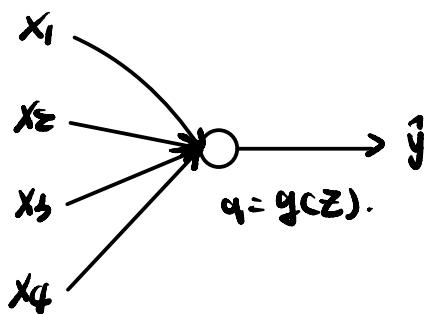
$\hat{y} = (w^{[L]})^T x \Rightarrow \hat{y}$ and da will explode.

if $w^{[l]} < I$, e.g. $w^{[L]} = \begin{bmatrix} 0.5 & 0 \\ 0 & 0.5 \end{bmatrix}$

$\hat{y} = (w^{[L]})^T x \Rightarrow \hat{y}$ and da will vanish.

Solution for exploding / Vanishing Problem.

Single neuron example:



$$z = w_1 x_1 + w_2 x_2 + w_3 x_3 + \dots + w_n x_n + b$$

ignored for now.

Intuition: Large $n \rightarrow$ Smaller w_i . So that z won't explode or vanish.

$$\text{Thus: } \text{Var}(w_i) = \frac{1}{n} \Rightarrow \text{Var}(n w_i) = \underline{1}$$

Prevent w_i from being too big or too small.

Implementation:

$$\text{Var}(W) \approx 1 \sim 2$$

$$W^{[l]} = np.random.randn(\text{Shape}) * \boxed{\sqrt{\frac{2}{n^{[l-1]}}}}$$

Mean at 0

$\frac{2}{n}$ rather than $\frac{1}{n}$

works best for ReLU.

Other variants:

tanh:

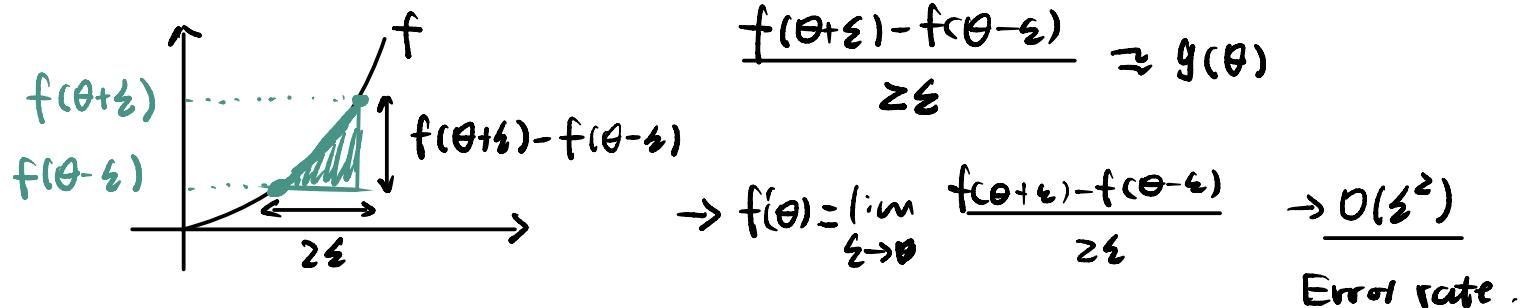
$$\sqrt{\frac{1}{n^{[l-1]}}}$$

Xavier Initialization:

$$\sqrt{\frac{2}{n^{[l-1]} + n^{[l]}}}$$

Gradient Checking

Numerical approximation of gradients



Gradient checking for a neural network:

- Take $W^{[1]}, b^{[1]}, \dots, W^{[L]}, b^{[L]}$ and reshape into a big vector θ .
- Take $dW^{[1]}, db^{[1]}, \dots, dW^{[L]}, db^{[L]}$ and reshape into a big vector $d\theta$.

Grad check:

for each i :

$$\rightarrow d\theta_{\text{approx}}^{[i]} = \frac{J(\theta_1, \theta_2, \dots, \theta_i + \epsilon, \dots) - J(\theta_1, \theta_2, \dots, \theta_i - \epsilon, \dots)}{2\epsilon}$$

$$\approx d\theta^{[i]} = \frac{\partial J}{\partial \theta^i} \quad | \quad d\theta_{\text{approx}} \stackrel{?}{\approx} d\theta.$$

$$\text{check } \frac{\|d\theta_{\text{approx}} - d\theta\|_2}{\|d\theta_{\text{approx}}\|_2 + \|d\theta\|_2} \approx \begin{cases} 10^{-7} & : \text{great} \\ 10^{-5} & : \text{OK} \\ 10^{-3} & : \text{worry} \end{cases}$$

$\epsilon = 10^{-7}$

Gradient Checking Implementation Notes:

- Don't use it in training - Only to debug
- if grad check fails, check components with biggest difference to identify bug.
- Remember the effect of regularization
- Doesn't work with dropout
- Run with random Initialization. Perhaps again after some training.